

Just the Two of Us? A Family of *Pseudomonas* Megaplasms Offers a Rare Glimpse into the Evolution of Large Mobile Elements

Brian A. Smith^{1,*}, Courtney Leligdon¹, and David A. Baltrus^{1,2}

¹School of Plant Sciences, University of Arizona

²School of Animal and Comparative Biomedical Sciences, University of Arizona

*Corresponding author: E-mail: basmith@email.arizona.edu.

Accepted: March 26, 2019

Data deposition: This project has been deposited at GenBank under the accessions NZ_CP032678.1, SRP157842, SRP157827, and CP031225-CP031228.

Abstract

Pseudomonads are ubiquitous group of environmental proteobacteria, well known for their roles in biogeochemical cycling, in the breakdown of xenobiotic materials, as plant growth promoters, and as pathogens of a variety of host organisms. We have previously identified a large megaplasmid present within one isolate of the plant pathogen *Pseudomonas syringae*, and here we report that a second member of this megaplasmid family is found within an environmental Pseudomonad isolate most closely related to *Pseudomonas putida*. Many of the shared genes are involved in critical cellular processes like replication, transcription, translation, and DNA repair. We argue that presence of these shared pathways sheds new light on discussions about the types of genes that undergo horizontal gene transfer (i.e., the complexity hypothesis) as well as the evolution of pangenomes. Furthermore, although both megaplasms display a high level of synteny, genes that are shared differ by over 50% on average at the amino acid level. This combination of conservation in gene order despite divergence in gene sequence suggests that this Pseudomonad megaplasmid family is relatively old, that gene order is under strong selection within this family, and that there are likely many more members of this megaplasmid family waiting to be found in nature.

Key words: megaplasmid, horizontal gene transfer, complexity hypothesis, pangenome, comparative genomics.

Introduction

Horizontal gene transfer (HGT) of megaplasms can rapidly create dramatic phenotypic differences between otherwise closely related bacterial strains, with potential for over a thousand genes to be gained by a strain in a single event. Although there have been numerous attempts to identify overarching themes for the evolutionary effects of HGT based on types of genes and pathways transferred, such efforts have often neglected to incorporate intrinsic characteristics of megaplasms (Lercher and Pál 2008; Wellner and Gophna 2008; Cohen et al. 2011; McInerney et al. 2017; Shapiro 2017; Vos and Eyre-Walker 2017). Furthermore, because secondary replicons are prone to rapid reshuffling of gene order as well as extensive gains and losses of loci, it has traditionally been challenging to analyze past evolutionary dynamics to understand overall historical pressures acting on this class of mobile

elements (Holden et al. 2004, 2009; Choudhary et al. 2007; Guo et al. 2009; Cooper et al. 2010; Janssen et al. 2010; Epstein et al. 2012; diCenzo and Finan 2017). More thorough investigation of evolutionary dynamics within relatively large plasmid families could therefore provide new viewpoints into the evolutionary effects of gene transfer and may also enable broader generalizations about selective forces driving the composition and overall structure of megaplasms, chromids, and secondary chromosomes.

Megaplasms are generally characterized as low copy extrachromosomal replicons >350 kb in size and which are dispensable to the bacterial cell under a subset of conditions (diCenzo and Finan 2017). As with many plasmid families, they have often been identified because they impart beneficial phenotypes such as resistance to antimicrobial compounds or introduce novel catabolic pathways into host cells

© The Author(s) 2019. Published by Oxford University Press on behalf of the Society for Molecular Biology and Evolution.

This is an Open Access article distributed under the terms of the Creative Commons Attribution Non-Commercial License (<http://creativecommons.org/licenses/by-nc/4.0/>), which permits non-commercial re-use, distribution, and reproduction in any medium, provided the original work is properly cited. For commercial re-use, please contact journals.permissions@oup.com

(diCenzo and Finan 2017). Given their size and gene content, it is possible that megaplasmiids possess greater potential for generating evolutionary costs than smaller plasmids when transferred to naive hosts (Baltrus 2013; San Millan and MacLean 2017). However, efforts to identify shared genes and pathways across megaplasmiids and to use this information to make predictions about potential systems level conflicts have been hampered by poor sampling across novel megaplasmiid families (diCenzo and Finan 2017). Identification of additional examples can help to fill this gap in current knowledge and may uncover new evolutionary trends that govern megaplasmiid–chromosomal interactions.

Consideration of megaplasmiids may uniquely inform general discussions about evolutionary effects HGT in ways that have been overlooked by analyses focusing simply on distributions of genes and pathways maintained after transfers and without considering timing of HGT or linkage between loci. For instance, many of the earliest discussions concerning evolutionary constraints of HGT, the so-called “complexity hypothesis,” found that loci associated with “more complex” cellular processes undergo lower rates of transfer than other genes (Jain et al. 1999; Cohen et al. 2011). Interpretations of these patterns have changed through time with examples of horizontally acquired informational genes, suggesting that it is actually the shape of protein interaction networks that are critical for maintenance after gene transfer (Tett et al. 2007; Cohen et al. 2011; Hall et al. 2015; Harrison et al. 2015; MacLean and San Millan 2015; Kacar et al. 2017). However, larger mobile elements like megaplasmiids can contain genes encoding proteins and pathways that could be classified as “complex” (i.e., proteins involved in translation) and which have the potential to interact with numerous chromosomally encoded pathways (Baltrus et al. 2011; diCenzo and Finan 2017). Likewise, a variety of recent articles have focused on selective forces (or lack thereof) governing microbial pangenomes (Andreani et al. 2017; McInerney et al. 2017; Shapiro 2017; Vos and Eyre-Walker 2017). These discussions have largely focused on population level distributions for single genes that compose a pangenome, but by their nature intrinsically fail to consider linkage of genes on plasmids that are frequently acquired and lost. Although many genes within the pangenome may indeed be “adaptive,” such a viewpoint overlooks the idea that no single gene need be adaptive for the bacterial cell if selection acts at the level of plasmid transfer and hundreds of genes that could be linked to that process. More thorough characterization of multiple megaplasmiid families and identification of new megaplasmiids will enable identification of the types of genes and pathways canonically associated with these large vectors and patterns that emerge can be incorporated into greater discussions of the general role of HGT across bacterial species.

Previously, we have described a megaplasmiid, pMPPla0107, found within one isolate of the phytopathogen *Pseudomonas syringae* (Baltrus et al. 2011). pMPPla0107 is self-

transmissible across the *Pseudomonas* phylogeny, harbors numerous loci that could be annotated as “housekeeping” genes, and it is stably maintained within recipient cells (Baltrus et al. 2011; Dougherty et al. 2014; Romanchuk et al. 2014). We have also demonstrated that acquisition of this megaplasmiid through HGT also imparts significant phenotypic costs to recipient cells, likely mediated by detrimental interactions between chromosomal and plasmid encoded proteins (Dougherty et al. 2014). It is unclear if pMPPla0107 is the only megaplasmiid of its kind and size with conjugation abilities across an entire genus of bacteria or if it is a member of a larger family of secondary replicons.

Here, we identify a new megaplasmiid, pBASL58, related to pMPPla0107 and use molecular and computational approaches to characterize this megaplasmiid family more broadly. We show that, although these megaplasmiids are similar in size, genetic structure, nucleotide bias, and functionality, there is a high level of divergence across shared orthologous gene groups, a dissimilar cargo region, and differing CRISPR loci. This overall level of divergence suggests that both members of this megaplasmiid family have been independently evolving for a relatively long period of time. Even more, we find that these divergent orthologous pathways demonstrate high levels of synteny in the context of overall plasmid structure; suggesting that conservation of gene orientation and order is under relatively strong selective pressures. Lastly, characterization of these plasmids allows for a chance to emphasize that pathways found on both megaplasmiids are likely involved in important cellular processes like nucleotide synthesis and DNA replication, which highlights new discussion points to add to the complexity hypothesis as well as the adaptive nature of pangenomes.

Materials and Methods

Identification of pBASL58

Initial BlastP searches coupled with inspection of contigs from a draft genome assembly of *Pseudomonas* sp. Leaf58 (Bai et al. 2015) suggested that this strain could contain a megaplasmiid related to pMPPla0107, represented within contig Ga0102293_111 within the draft genome assembly containing 18 contigs total (GenBank accession GCA_001422615.1). An isolate of *Pseudomonas* sp. Leaf58 was obtained from DSMZ (DSM-102683), and a single colony was picked from a culture from the freeze-dried ampule plated on unsupplemented lysogeny broth media. To test for circularization of contig Ga0102293_111, primers (BAS 17-31) were designed to amplify off of the edges of the contig and overlap each other. An approximate 1.5-kb size polymerase chain reaction (PCR) product was amplified from an overnight culture of this strain grown in King’s B (KB) media, demonstrating circularization of this contig. Sanger sequencing of this PCR fragment demonstrated that the initial draft contig contained a missassembly, and so this sequence was corrected by hand, and the

contig was reoriented and used in all analyses in this article. Sequence of this contig can be found at Figshare (doi:10.6084/m9.figshare.6914033). For consistency of analyses throughout the article, we reannotated the megaplasmiid sequence with Prokka v1.12 using default parameters, and this annotation can also be found at Figshare (doi:10.6084/m9.figshare.6914033).

Genome Sequencing, Assembly, and Annotations of *Pseudomonas* sp. Leaf58

As further evidence of the existence of a megaplasmiid in Leaf58, we generated a complete genome assembly for this strain (currently found at Figshare [doi:10.6084/m9.figshare.6914033], GenBank accession TBD). After revival from the Baltrus lab stock, a single colony *Pseudomonas* sp. Leaf58 was picked to an overnight culture in KB media and grown in a shaking incubator at 27 °C. After ~24 h, DNA was extracted from this culture using a Promega Wizard kit. A rapid sequencing library was created using this DNA, and 169,316 reads (933,937,907 total bp, 5,515 bp average read size) were generated on an Oxford Nanopore MinION R9.4 flowcell using a Rapid sequencing kit (SQK-RAD004) and can be found in the Sequence Read Archive (SRA) (accession SRP157842). Additionally, 100-bp paired end Illumina reads used to generate the original draft genome of this strain were downloaded from the SRA (accession ERR1103815) (Bai et al. 2015). A complete genome sequence for *Pseudomonas* sp. Leaf58 was generated by combining these short and long reads in Unicycler (version 0.4.4) (Wick et al. 2017). This sequence consists of a single chromosome (5,432,868 bp) and the pBASL58 megaplasmiid (904,253 bp), both of which were circular according to Unicycler.

Genome Sequencing, Assembly, and Annotations of Pla107

A single colony of the Baltrus lab stock of *P. syringae* pv. *lachrymans* 107 (MAFF31015) was picked to an overnight culture in KB media, and grown in a shaking incubator at 27 °C. After ~24 h, DNA was extracted from this culture using a Promega Wizard kit. Illumina sequencing of Pla107 was performed by MicrobesNG, and generated 2,771,213 250-bp paired end reads (231 median read length after trimming, ~166× coverage of the genome) on an Illumina MiSeq. Assembly was performed using SPAdes v3.10.1 with default parameters as well as through MicrobeNG's bioinformatics pipeline, which matches the reads to the best reference using Kraken and maps reads back to that reference using BWA-MEM (Wood and Salzberg 2014). MicrobeNG also uses de novo assembly with SPAdes. pMPPla107 assembled completely from these reads, and this version of the megaplasmiid sequence can be found Figshare (doi:10.6084/m9.figshare.6914033) and was used for all analyses throughout this article. Gene annotation of this version of the

megaplasmiid sequence was performed with Prokka v1.12 using default parameters. This gene model used for all coding sequence analyses within the article and can be found at Figshare (doi:10.6084/m9.figshare.6914033). We additionally generated long read sequences for Pla107 using a MinION from Oxford Nanopore. A rapid sequencing library was created from an independent genomic isolation of a derivative of Pla107, DBL328, which contains an integrated version of the pMTN1907 marker plasmid and which has been selected to for kanamycin resistance from this marker plasmid. As above, a single colony of this strain was picked to an overnight culture in KB media and DNA was extracted with a Promega Wizard kit. Fifteen thousand four hundred sixty-one reads (139,041, 576 total bp, 8,993 average read size) were generated on an R9.4 flowcell using a Rapid sequencing kit (SQK-RAD004). A whole genome assembly was created by combining both MiSeq and MinION reads using Unicycler (version 0.4.4) (Wick et al. 2017) with default parameters. Reads from the Illumina and Oxford Nanopore platforms can be found in the SRA (accession SRP157827). This whole genome sequence consists of a circular chromosome (6,075,120 bp), pMPPla107 (971,889 bp, and sequence identical to the assembly from SPAdes alone), and two other plasmids, pPla107-1 (62,136 bp) and pPla107-2 (40,720 bp). Three of these sequences (except pPla107-1) were complete and circular contigs according to Unicycler assembly. This assembly was used to update the GenBank version of this genome and is found at accessions (CP031225, CP031226 CP031227, and CP031228). Gene annotations in this GenBank file were generated by National Center for Biotechnology Information's PGAAP pipeline (Tatusova et al. 2016).

Identifying Origins of Replication

To identify putative origins of replication for both megaplasmiids, we used a modified Guanine and Cytosine (GC) skew script (Charif and Lobry 2007) to scan the entirety of pMPPla107 and pBASL58 and combined this information with characterization of repetitive motifs that could represent *oriV* sites. GC skew and repetitive motifs suggest pMPPla107 and pBASL58 have predicted origins of replication within a similar genomic region near partitioning genes (fig. 1). Based on this information we oriented the sequences of pMPPla107 and pBASL58 to begin at the start codon of shared *parA*-like loci. We chose the *parA*-like locus as the starting point because it is shared by both sequences, is near the predicted origin of replication, and is predicted to be an important gene for plasmid partitioning.

CRISPR Identification

CRISPR-Cas and repeat structure annotations were identified using both Prokka annotations and the web tool CRISPRCasFinder (Grissa et al. 2007; Abby et al. 2014; Couvin et al. 2018).

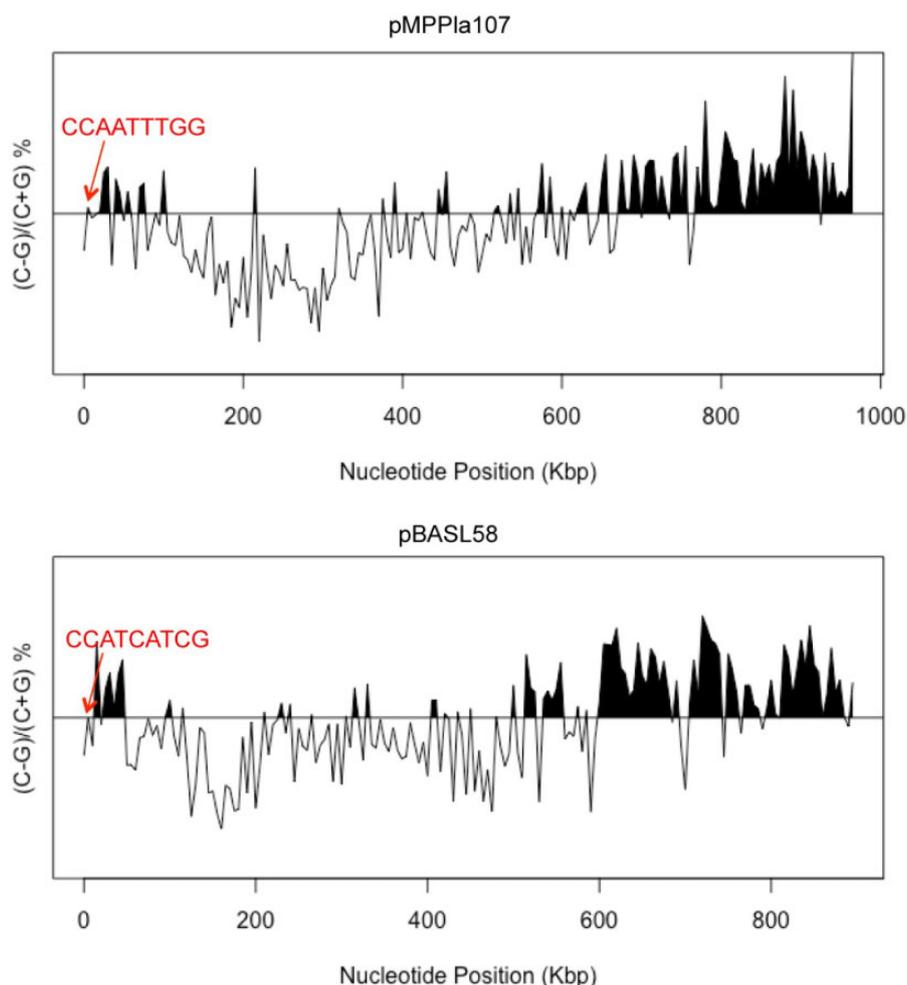


Fig. 1.—Predicted origins of replication occur within a similar region of pMPPla107 and pBASL58. GC skew was calculated using a sliding window of 5 kb and is a known predictor of origins of replication indicated by a dramatic shift in GC content. Repetitive motifs were also calculated for areas near the predicted origin of replication as repetitive binding regions occur near replication sites. The most common motif is indicated in red.

Plasmid Comparisons with BlastP SynMap, and MAUVE

Amino acid sequence names were changed to numbers in an increasing order using the `mod_protein_id.py` script. We then used the BLAST 2.6.0+ package (Altschul et al. 1990). BlastP parameters were altered to ensure only the top hit was returned and that there were zero overlapping hits. The BLAST command used was

```
blastp -db [blastdb] -query [query_file]
  -culling_limit 1 -max_target_seqs 1 -max_hsp 1
  -out [out_file] -outfmt 6
```

Data were extracted from the BLAST output at 40%, 50%, 60%, and 70% identity cutoffs and plotted in R using `ggplot2`.

Sequences for pMPPla107 and pBASL58 were uploaded and input to the CoGe web software SynMap using the Last algorithm and default parameters (Haug-Baltzell et al. 2017).

pMPPla107 and pBASL58 sequences were input into Progressive Mauve 2.3.1 to compare megaplasmid sequences within the software Geneious (Darling et al. 2010).

Gene Mapping Visualization with Circa

The BlastP output data mentioned above was altered in a format to comply with input to Circa using `gff_info_extract.py` followed by `geneid_match.py`. The parameters used to generate the Circa map and the Python scripts used to generate the data can be found at the https://github.com/basmith89/megaplasmid_compare; last accessed February 5, 2019.

Tetranucleotide Frequency Comparisons

We performed pairwise comparisons of tetranucleotide frequencies between chromosome sequences and secondary replicon sequences in an all by all method.

Tetranucleotide frequencies were calculated with the `calc.kmerfreq.pl` script created by Mads Albertsen (Albertsen et al. 2013) found at <https://github.com/MadsAlbertsen/multi-metagenome>; last accessed February 5, 2019. Output of this script was plotted using `ggplot2` and R^2 values were calculated in R.

Functional Comparisons with KEGG, KASS, and UProC

We carried out two analyses utilizing the Kyoto Encyclopedia of Genes and Genomes (KEGG) database (Kanehisa and Goto 2000). Amino acid sequences of coding regions predicted by Prokka were input into the protein sequence classification software, UProC (Meinicke 2015). UProC's output is a list of KEGG IDs and counts. We designed a perl script, `kegg_path_counter.pl`, to extract these ID's and counts and associated them with KEGG functional pathways. The script and ID key can be found at https://github.com/basmith89/megaplasmid_compare. These data were then plotted with the `ggplot2` package in R.

Amino acid sequences output by Prokka were also run through a Python script to produce a list of gene annotations that both megaplasms have in common https://github.com/basmith89/megaplasmid_compare. Amino acid sequences from genes on this shared list were then run through KASS (KEGG Automatic Annotation Server) to determine what pathways are shared by the megaplasms (Moriya et al. 2007). Pathways maps were then condensed into one figure by hand.

Results

A New Member of the pMPPla107 Megaplasms Family

pMPPla107 was originally identified from an assembly using both 454- and 30-bp Illumina sequencing reads (Baltrus et al. 2011). However, due to limitations of these early technologies, this assembly of pMPPla107 remained incomplete and consisted of linked scaffolds. We therefore utilized updated sequencing and assembly technologies to sequence the *P. syringae* genome containing pMPPla107, yielding a complete circular sequence for this megaplasms (971,889 bp compared with 963,598 bp in original sequence) (table 1). Additionally, multiple searches using protein sequences from pMPPla107 consistently yielded high quality matches to the scaffold Ga0102293_111 (referred to as pBASL58 from here on) from a public genome assembly of *Pseudomonas* sp. Leaf58. This strain was originally isolated as part of a project to thoroughly sample cultureable strains from the phyllosphere of Arabidopsis and is most closely related to *Pseudomonas putida* strains (Hesse et al. 2018). We independently confirmed circularization (fig. 2) of this contig from Leaf 58, using both PCR and long read nanopore sequencing, definitively showing this contig was indeed a large megaplasms separate from the chromosome.

Table 1

General Features of *Pseudomonas syringae* and *Pseudomonas* Leaf58 Replicons Have Similarities

Name	Size	Genes (CDS)	tRNA	GC Content
pMPPla107	971,871	1,082	54	52.84
pBASL58	903,765	996	44	55.4
Leaf58 Concatenated	5,378,738	4,847	80	62.35
Lac107 Concatenated	5,936,302	5,436	58	58.26

NOTE.—Size, coding regions, tRNAs, and GC content were calculated to understand the relationship between the four sequences on a broad scale. "Concatenated" indicates that all sequences from the assemblies, but disregarding either the pMPPla107 or pBASL58 replicons to their respective genome were concatenated together.

Both Megaplasms Contain Numerous tRNA Loci

The size, number of predicted genes, number of tRNAs, and GC content are highly similar between pMPPla107 and pBASL58 (table 1). Overall GC content was similar in Leaf58 and *P. syringae lac107*, and the GC content in both pMPPla107 and pBASL58 were lower than their respective chromosomal partners. pBASL58 and pMPPla107 contained 54 and 44 regions annotated as tRNA loci, respectively. pBASL58 encodes 20 unique tRNAs and pMPPla107 encodes 10, some of which were repetitive like tRNA-Glu(ttc) in pBASL58 occurring six times. When observing tRNA amino acid products, pMPPla107 encodes for 16/20 possible amino acids and pBASL58 encodes for 19/20 possible amino acids possibly indicating pBASL58 is less dependent on host tRNAs. In addition to the 16 amino acids produced by pMPPla107, pBASL58 is predicted to code for the ability to charge tRNAs with tryptophan, glutamate, and aspartate and both plasmids are missing any anticodons to produce histidine. These differences could suggest an amino acid preference for the maintenance or protein production of the plasmids.

Identifying Genomic Similarities of pMPPla107 and the Leaf58 Plasmid

Both megaplasms within this new family are highly syntenic (fig. 3A), with Mauve alignment showing that 72.8% (707,677 bp out of 971,871 bp) of pMPPla107 aligns well with 71.6% (646,763 bp out of 903,765 bp) of pBASL58 (supplementary fig. 1, [Supplementary Material](#) online). The regions of highest similarity occur near the origin of replication. Despite overall high levels of synteny, there is a highly dissimilar region (~300 kb in size) occurring within the first half of the sequences and a ~50-kb inversion in the last half indicating these megaplasms have also undergone structural diversification.

Even though both megaplasms display high levels of synteny, preliminary comparisons of protein sequences suggested a relatively high level of divergence between orthologs shared by both megaplasms (fig. 3B and C). The highest levels of average amino acid similarity (48.6%) occur near

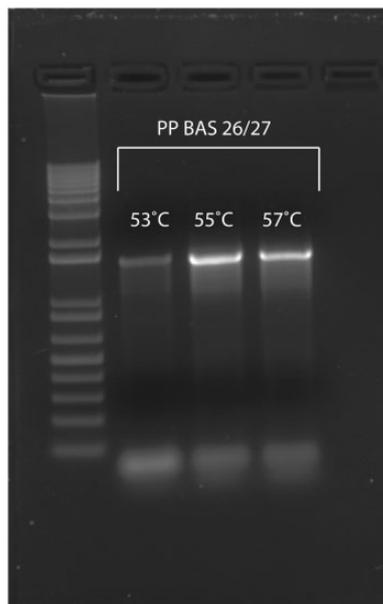


FIG. 2.—Confirmation of circular DNA molecule of Leaf58 megaplasmid. Primers designed to amplify the ends of the Leaf58 contig and disregarding the misassembled repeat region successfully amplified products of an expected size. Several primer sets were used and this image is one primer set that spanned the overlapping region. To confirm that PCR products from the correct site, all samples underwent Sanger sequencing and were mapped back to the genome sequence. Three annealing temperatures (53, 55, and 57°C) were used due to difficulties amplifying this region.

the predicted origin of replication where genes for plasmid replication, partitioning, and conjugation are common. Areas near the terminus still demonstrate strong synteny but have higher divergence in amino acid identity ($\approx 38.2\%$ similarity). These data suggest pMPPla107 and pBASL58 are structurally related to each other and share a common plasmid ancestor but have experienced independent evolutionary pressures for long enough time for significant diversification to occur within shared protein sequences.

To further gauge relationships between both megaplasmids and the chromosomes of their host strains, we compared tetranucleotide frequencies for each of these replicons (Teeling et al. 2004; Richter and Rosselló-Móra 2009; Nishida et al. 2012). Pairwise comparisons demonstrated that pMPPla107/pBASL58 ($R^2 = 0.878$) and the *P. syringae*/Leaf58 ($R^2 = 0.889$) chromosomes are most similar in frequencies (fig. 4). All remaining pairwise comparisons reported R^2 values less than 0.780. pMPPla107 shows the greatest differences in tetranucleotide frequencies when compared with both the *P. syringae* and the Leaf58 chromosomes with R^2 values of 0.524 and 0.393, respectively. pBASL58 shares slightly more similar frequency preferences indicative of R^2 values of 0.780 and 0.695, to *P. syringae* and Leaf58 chromosomes, respectively. These data suggests that mutational biases affecting these secondary replicons are most

similar to each other, which suggests that they have not been replicating within these host strains long enough to be subject to amelioration.

Housekeeping Gene Functionality Is Shared by pBASL58 and pMPPla107

Based on the structural similarities established, we hypothesized that pMPPla107 and pBASL58 would share similar functional pathways. UProC called 9% (85) and 10% (111) of the predicted coding regions for pBASL58 and pMPPla107, respectively, indicating the majority of predicted gene functionality is unknown. Annotation with Prokka returned similar results (13% of genes with annotated functions). The pathways and functions most frequently annotated were replication and repair at 2.3% (22 genes) for pBASL58 and 2.4% (26) for pMPPla107, global and overview maps at 2.1% (20) for pBASL58 and 2.2% (24) times for pMPPla107, and nucleotide metabolism at 1.3% (12) for pBASL58 and 1.8% (19) for pMPPla107 (fig. 5). KEGG KASS also predicted that the two megaplasmids share 57.6% (99/172) of annotated genes. Therefore, pBASL58 and pMPPla107 carry 31 and 42 unique genes, respectively. Again, the overall distribution of gene products present on both megaplasmids tends toward DNA synthesis, DNA repair, and synthesis of deoxyribonucleotide-triphosphates (supplementary table 1 and fig. 3, [Supplementary Material](#) online). These shared groups include DNA polymerase III subunits, helicases, primase, ligases, recombination proteins, and exonucleases indicating these megaplasmids encode for pathways associated with their maintenance. Other gene products on these megaplasmids are involved in metabolic pathways such as fatty acid biosynthesis, RNA degradation, Aminoacyl tRNA biosynthesis, and NOD-like receptor signaling pathways. Interestingly, both plasmids also encode for several membrane and multidrug efflux pump genes. Both shared efflux genes belong to the Resistance–Nodulation–Division family of transporters and are known for their multidrug resistance efflux capabilities indicating potential selective factors enabling maintenance in host cells.

Differences of pMPPla107 and the Leaf58 Plasmid

pBASL58 is predicted to encode a complete CRISPR system from 229 to 241 kb, including two *cas*, three *csy* genes, and a repeat region that includes 36 repeats and spacers (fig. 6). This CRISPR is located in the region of dissimilarity between pMPPla107 and pBASL58 and is not found in pMPPla107. pBASL58 and pMPPla107 do share an (presumably) incomplete CRISPR systems at 436 and 576 kb, respectively (fig. 6). These regions include *cas3*, *csy3*, and *csy4* but lack *csy1* and *csy2*. pMPPla107 lacks a repeat region altogether associated with this locus, whereas pBASL58 has a repeat region at 720 kb encoding nine repeats and spacers. To our knowledge,

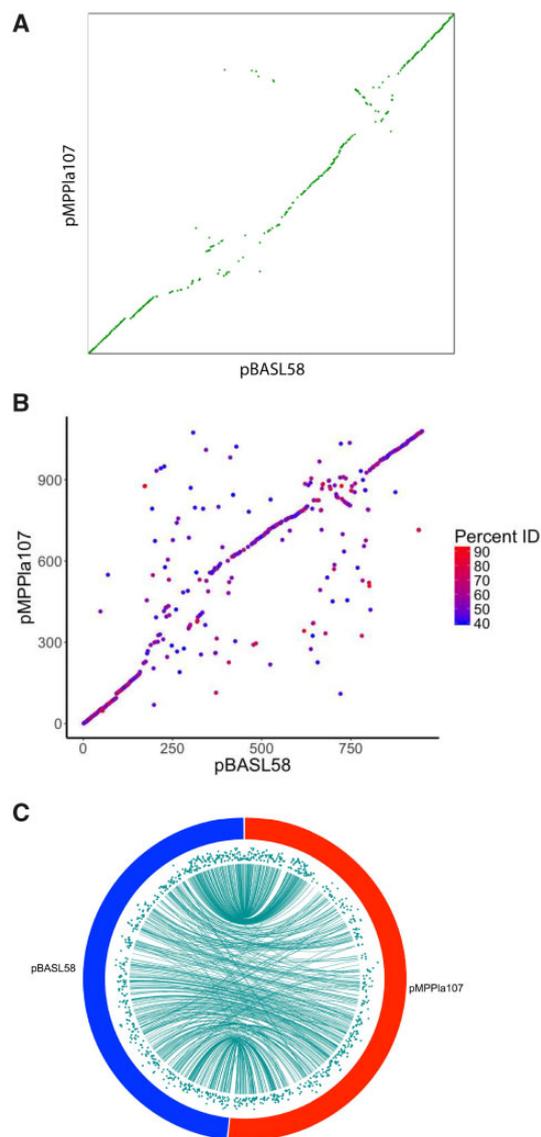


FIG. 3.—pMPPla107 and pBASL58 share synteny and demonstrate divergence on the amino acid level. (A) SynMap output of pMPPla107 versus pBASL58 sequences suggests highly syntenic megaplasms. x axis is pBASL58 gene order where $x_{1...N} = \text{gene}_{1...N}$, and the y axis is pMPPla107 gene order where $y_{1...N} = \text{gene}_{1...N}$. Completely syntenic sequences would be represented by $y = 1x + b$. (B, C) BLAST data was used to plot pMPPla107 versus pBASL58 syntenicity and amino acid divergence data together. (B) The majority of syntenic genes have $\geq 50\%$ sequence identity. x axis is pBASL58 gene order where $x_{1...N} = \text{gene}_{1...N}$, and the y axis is pMPPla107 gene order where $y_{1...N} = \text{gene}_{1...N}$. Completely syntenic sequences would be represented by $y = 1x + b$. The color gradient is set to BlastP percent identity results for each comparison. A percent identity cutoff of $\geq 40\%$ was used. (C) Circa plot using BLAST data indicates higher synteny near the origin, whereas areas near the terminus are less syntenic and experience more inconsistent mapping. Teal lines connect gene start position on pBASL58 to gene start position on pMPPla107. Teal scatter plots are amino acid sequence identity with $40\% = 0$ (bottom) and $100\% = 100$ (top).

these are the first complete CRISPR systems located on plasmids found within Proteobacteria.

There exists a region of dissimilarity across both megaplasms, occurring after ~ 170 kb (fig. 3), which could be classified as a cargo region. In pMPPla107 this region consists of 468 predicted genes, of which 27 are annotated. Eighteen of these 27 annotated genes can be found in pBASL58 and again encode for genes associated with DNA replication, repair, and metabolism. These genes also include membranous proteins like FtsH, which is known to degrade unnecessary or damaged membrane proteins (Tomoyasu et al. 1995; Ito and Akiyama 2005). We have also found that this region can largely be deleted from pMPPla107 during lab adaptation (unpublished) even though the rest of the plasmid is maintained. These data suggest that although this large region may be expendable in some strains, pBASL58 has maintained many of the annotated genes perhaps pointing to their importance in megaplasmid stability or maintenance.

Discussion

We report a family of divergent, yet syntenic megaplasms found in single isolates across distinct *Pseudomonas* species. High levels of synteny are matched by shared signals in both tetranucleotide bias and protein pathway functionality. However, these plasmids hosted by strains that are phylogenetically and geographically separated; *Pla107* (containing pMPPla107) was found within a *P. syringae* isolate as a causative agent of cucumber disease in Japan, whereas *Leaf58* was found as an epiphyte of *Arabidopsis* in Switzerland in a strain most closely related to *Pseudomonas putida* (Hesse et al. 2018). Furthermore, despite high levels of synteny and shared protein functionality, consistently high levels of divergence across shared proteins ($\approx 30\%$) suggest both plasmids have been independently evolving for a relatively long period of time. From these data, we infer that multiple additional members of a family of relatively large (≈ 1 Mb) “cryptic” megaplasms likely persist within *Pseudomonas* strains.

That there have been no signs of these megaplasms in the numerous sequences of *Pseudomonads* closely related to each of these isolates is strong indication that these megaplasms have been relatively recently acquired by their host strains. This pattern, coupled with high levels of divergence between members of this megaplasmid family, suggest that these replicons likely have a high turnover rate within strains over evolutionary time and may persist within communities through frequent horizontal transfer. In other words, presence of this megaplasmid family may be transient in any given genome, but has likely been maintained within *Pseudomonads* for a long time. Such a lifestyle is consistent with high levels of conjugation as observed in pMPPla107 under laboratory conditions (Romanchuk et al. 2014).

Replication, transcription, and translation of horizontally transferred genes are known to incur costs on host cell

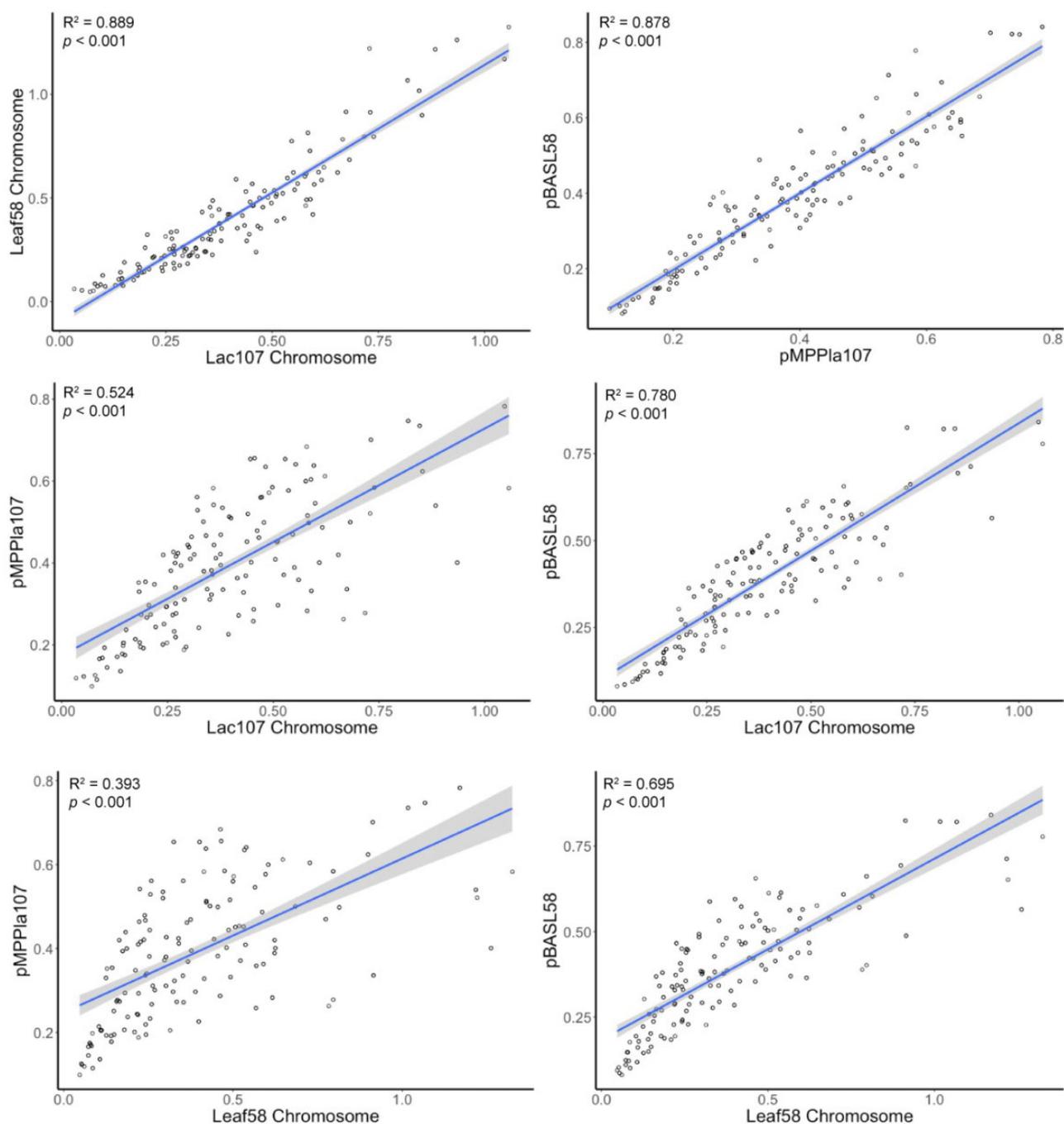


FIG. 4.—Tetranucleotide frequencies between pMPPla107 and pBASL58 suggest an evolutionary relationship. Nucleotide biases were determined to demonstrate relatedness of megaplasmid and chromosomal sequences in a pairwise fashion. The blue line represents the linear regression model with the surrounding shaded gray area indicating a 95% confidence interval. R^2 and P values are listed for each comparison.

resources with protein production likely having the greatest effect on fitness (Bragg and Wagner 2009; Shachrai et al. 2010; Baltrus 2013; Hall et al. 2017). Previous work on pMPPla107 suggests that acquisition of the megaplasmid results in lowered fitness and other phenotypic changes which could be costly in some environments, yet it still transfers readily and is maintained within host cells

(Dougherty et al. 2014; Romanchuk et al. 2014). Such costs could likely be the reason pMPPla107 and pBASL58 encode a large number of genes involved in critical functions regarding plasmid maintenance and transmission as well as potential addiction systems and could enable long-term survival despite a transient lifestyle. In particular, there are various proteins found in pMPPla107 and pBASL58 involved in synthesizing

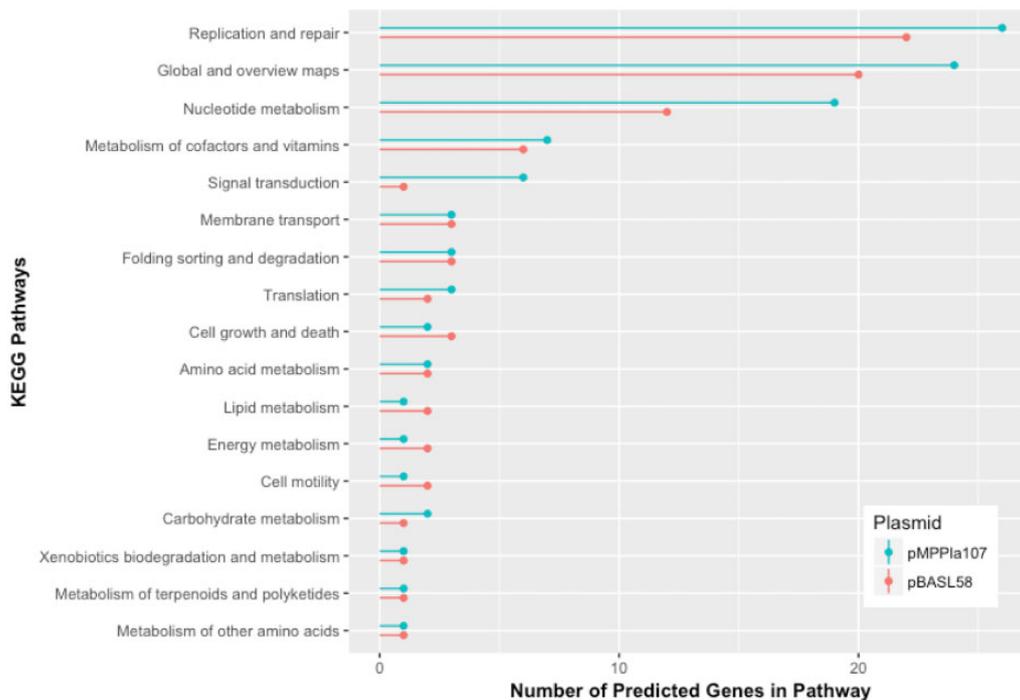


FIG. 5.—pMPPla107 and pBASL58 share similar functional profiles. pMPPla107 and pBASL58 were input into UProC which counts the number of gene sequences that are predicted to belong to a KEGG ID. KEGG IDs were then matched with the correct pathway. Counts account for only a small subset of predicted genes on both megaplasmids, as most predictions are hypothetical proteins and could not be assigned a KEGG pathway.

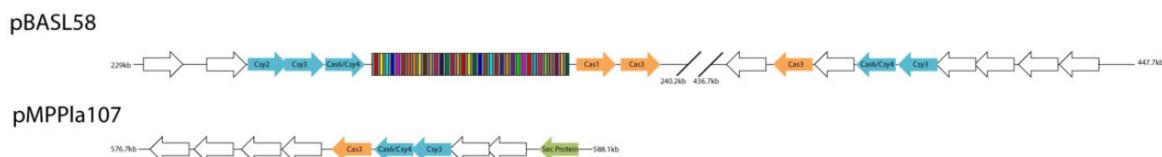


FIG. 6.—CRISPR systems on pBASL58 and pMPPla107. pBASL58 encodes two CRISPR loci, one of which contains a repeat-spacer regions of 36 repeats. pMPPla107 contains a CRISPR locus without any repeat-spacer regions. Direction of arrows indicates gene orientation. Arrows are colored as (blue) *Cys* genes, (orange) *Cas* genes, (green) secretion genes, and (white) hypothetical genes. The multicolored boxes indicate the repeat-spacer region, where gray boxes are spacers and colored boxes are repeats.

precursors for nucleotides such as thymidylate synthase, guanylate kinase, ribonucleoside diphosphate reductase, deoxycytidine triphosphate deaminase, and glutamate synthase (supplementary table 1 and fig. 3, [Supplementary Material](#) online). The megaplasmids may carry these proteins in order to increase flux to nucleotide synthesis and drive replication and transcription processes to alleviate any physiological costs an additional ≈ 1 Mb of newly acquired DNA may bring. Many of these genes do not encode for complete pathways, indicating possible parasitic behavior of host resources while ensuring the necessary building blocks for plasmid maintenance are available.

Plasmid usage of host tRNA pools has been shown to deplete tRNAs resulting in reduced growth and fitness (Elf et al. 2003; Bonomo and Gill 2005; Dittmar et al. 2005; Harrison et al. 2015). The large number of tRNAs and presence of a

handful of annotated tRNA ligases encoded on the megaplasmids may serve the purpose of avoiding translational costs due to tRNA depletion or may accommodate codon usage bias between chromosome and megaplasmid. Both megaplasmids are also predicted to encode Mfd, Rep, DnaB, and RecA all known to resolve replication and transcription complex conflicts ensuring successful replicon duplication and transcription (McGlynn et al. 2012; Hamperl and Cimprich 2016). We hypothesize the megaplasmids maximize their ability to persist by eliminating or compensating for these potential costs by encoding a variety of housekeeping genes coupled with high levels of horizontal transfer through conjugation.

Evolutionary relationships between pBASL58 and pMPPla107, their relatively large size and contribution to gene content of single strains, coupled with maintenance of

“housekeeping” genes, and high levels of transfer across Pseudomonads suggest that this megaplasmid will provide unique insights into an evolutionary argument concerning horizontal transfer referred to as the complexity hypothesis (Jain et al. 1999). The complexity hypothesis has been through multiple revisions but is currently interpreted as a trend where horizontally transferred genes are less likely to be involved with complex processes (like translation) and maintain a lower number of protein–protein interactions than vertically inherited loci (Cohen et al. 2011). One current limitation of the complexity hypothesis, as highlighted by these megaplasmid families, is that it fails to reconcile gene conservation in the context of highly mobile selfish DNA like plasmids. Both pBASL58 and pMPPla107 contain numerous “complex” genes, including those involved in nucleotide synthesis, DNA replication, and translation and yet these genes are clearly horizontally transferred across strains. Therefore, the presented family of megaplasmids potentially necessitates a caveat to the complexity hypothesis in which “complex” genes can be horizontally transferred frequently but are not maintained over time, because they are linked together on megaplasmids that require these pathways to ameliorate physiological costs.

Likewise, there have been numerous recent discussions about whether bacterial pangenomes are adaptive or neutral. Similar to the complexity hypothesis, these discussions tend to focus on the presence/absence of single genes across a variety of closely related genomes rather than the linked gain/loss of genes that compose a pangenome (McInerney et al. 2017; Shapiro 2017; Vos and Eyre-Walker 2017). To put this in perspective, recent findings suggest that the *P. syringae* pangenome is composed of 77,728 genes, meaning that 1.5% of these are solely present on pMPPla107 (Dillion et al. 2019). Because megaplasmids have the potential to add thousands of genes to a pangenome linked together in a single transfer event (Nowell et al. 2014), one has to consider that evolutionary pressures may act differentially on subsets of the pangenome. Our data suggest that a majority of genes on these megaplasmids may be either neutral or costly to the host when selection is considered in the context of the host genome. However, a majority of genes linked on the megaplasmid may be selectively beneficial for megaplasmid maintenance and/or transfer regardless of fitness of the host cell. Thus, presence of a majority of genes on the megaplasmid (and which are part of the pangenome) are under selection at some level, but only a minority of these may be beneficial at the level of bacterial strains or populations.

CRISPR-Cas systems have become popularized recently because of their utility in genome editing, however, these systems likely originated in bacteria as defense mechanisms against invasion of foreign genetic material (Deveau et al. 2010; Marraffini and Sontheimer 2010; Doudna and Charpentier 2014; Hsu et al. 2014; Makarova et al. 2015). CRISPR arrays are often carried and transferred by larger

plasmids in bacteria and archaea, yet *cas* genes are rarely found on plasmids (Godde and Bickerton 2006; Lillestøl et al. 2009). Here, we characterize a potentially shared CRISPR-Cas system bound to the bacterial megaplasmids pMPPla107 and pBASL58. Although pBASL58 encodes a fully intact CRISPR-Cas3 system with a region containing 36 spacers and repeats, this repeat and spacer region are not present within pMPPla107 leading us to believe pMPPla107’s system is nonfunctional. Regardless of functionality, it is quite interesting that at least one of these megaplasmids contains an intact CRISPR locus given the widespread idea that these systems are used by bacteria to defend against parasites and mobile elements. Perhaps the presence of a CRISPR system is a beneficial and selective trait for retention of pBASL58 in host cells in that it provides a transferable immune pathway. However, the recent description of CRISPR spacers that target sites on bacterial chromosomes also suggest that these loci may also function in gene regulation (Stern et al. 2010; Vercoe et al. 2013; Briner et al. 2015; Guan et al. 2017).

Using comparative computational and molecular approaches, we have characterized pBASL58, the second member of a family of large megaplasmids found in Pseudomonads. Conservation of pathway presence and megaplasmid structure strongly suggests that a majority of the sequences on pBASL58 and pMPPla107 have diverged from a common ancestral plasmid. However, the consistent levels of divergence between proteins shared by both plasmids suggest that this common ancestral plasmid did not recently exist. Finding two related plasmids with such high level of divergence also highlights the likelihood that other members of this megaplasmid family exist in nature and are waiting to be found. Our work serves as a guide to discover megaplasmid families as well as a foundation of understanding the forces that structure megaplasmid evolution, maintenance, and transfer.

Acknowledgments

This work was supported by internal funds from the University and US Department of Agriculture (USDA) NIFA 2016-67014-24805.

Supplementary Material

Supplementary data are available at *Genome Biology and Evolution* online.

Literature Cited

- Abby SS, Néron B, Ménager H, Touchon M, Rocha E. 2014. MacSyFinder: a program to mine genomes for molecular systems with an application to CRISPR-Cas systems. *PLoS One* 9(10):e110726.
- Albertsen M, et al. 2013. Genome sequences of rare, uncultured bacteria obtained by differential coverage binning of multiple metagenomes. *Nat Biotechnol.* 31(6):533–538.

- Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. 1990. Basic local alignment search tool. *J Mol Biol.* 215(3):403–410.
- Andreani NA, Hesse E, Vos M. 2017. Prokaryote genome fluidity is dependent on effective population size. *ISME J.* 11(7):1719–1721.
- Bai Y, et al. 2015. Functional overlap of the *Arabidopsis* leaf and root microbiota. *Nature* 528(7582):364–369.
- Baltrus DA. 2013. Exploring the costs of horizontal gene transfer. *Trends Ecol Evol.* 28:489–495.
- Baltrus DA, et al. 2011. Dynamic evolution of pathogenicity revealed by sequencing and comparative genomics of 19 *Pseudomonas syringae* isolates. *PLoS Pathog.* 7(7):e1002132.
- Bonomo J, Gill RT. 2005. Amino acid content of recombinant proteins influences the metabolic burden response. *Biotechnol Bioeng.* 90(1):116–126.
- Bragg JG, Wagner A. 2009. Protein material costs: single atoms can make an evolutionary difference. *Trends Genet.* 25(1):5–8.
- Briner AE, et al. 2015. Occurrence and diversity of CRISPR-Cas systems in the genus *Bifidobacterium*. *PLoS One* 10(7):e0133661.
- Charif D, Lobry JR. 2007. SeqinR 1.0-2: a contributed package to the R project for statistical computing devoted to biological sequences retrieval and analysis. In: *Structural approaches to sequence evolution. Biological and medical physics, biomedical engineering.* Berlin (Germany): Springer. p. 207–232.
- Choudhary M, Zanhua X, Fu YX, Kaplan S. 2007. Genome analyses of three strains of *Rhodobacter sphaeroides*: evidence of rapid evolution of chromosome II. *J Bacteriol.* 189(5):1914–1921.
- Cohen O, Gophna U, Pupko T. 2011. The complexity hypothesis revisited: connectivity rather than function constitutes a barrier to horizontal gene transfer. *Mol Biol Evol.* 28(4):1481–1489.
- Cooper VS, Vohr SH, Wrocklage SC, Hatcher PJ. 2010. Why genes evolve faster on secondary chromosomes in bacteria. *PLoS Comput Biol.* 6(4):e1000732.
- Couvin D, et al. 2018. CRISPRCasFinder, an update of CRISPRFinder, includes a portable version, enhanced performance and integrates search for Cas proteins. *Nucleic Acids Res.* 99: 7536.
- Darling AE, Mau B, Perna NT. 2010. progressiveMauve: multiple genome alignment with gene gain, loss and rearrangement. *PLoS One* 5(6):e11147.
- Deveau H, Garneau JE, Moineau S. 2010. CRISPR/Cas system and its role in phage-bacteria interactions. *Annu Rev Microbiol.* 64(1):475–493.
- diCenzo GC, Finan TM. 2017. The divided bacterial genome: structure, function, and evolution. *Microbiol Mol Biol Rev.* 81:e00019–17.
- Dillion MM, et al. 2019. Recombination of ecologically and evolutionarily significant loci maintains genetic cohesion in the *Pseudomonas syringae* species complex. *Genome Biol.* 20:3. doi: 10.1186/s13059-018-1606-y.
- Dittmar KA, Sørensen MA, Elf J, Ehrenberg M, Pan T. 2005. Selective charging of tRNA isoacceptors induced by amino-acid starvation. *EMBO Rep.* 6(2):151–157.
- Doudna JA, Charpentier E. 2014. Genome editing. The new frontier of genome engineering with CRISPR-Cas9. *Science* 346(6213):1258096.
- Dougherty K, et al. 2014. Multiple phenotypic changes associated with large-scale horizontal gene transfer. *PLoS One* 9(7):e102170.
- Elf J, Nilsson D, Tenson T, Ehrenberg M. 2003. Selective charging of tRNA isoacceptors explains patterns of codon usage. *Science* 300(5626):1718–1722.
- Epstein B, et al. 2012. Population genomics of the facultatively mutualistic bacteria *Sinorhizobium meliloti* and *S. medicae*. *PLoS Genet.* 8(8):e1002868.
- Godde JS, Bickerton A. 2006. The repetitive DNA elements called CRISPRs and their associated genes: evidence of horizontal transfer among prokaryotes. *J Mol Evol.* 62(6):718–729.
- Grissa I, Vergnaud G, Pourcel C. 2007. CRISPRFinder: a web tool to identify clustered regularly interspaced short palindromic repeats. *Nucleic Acids Res.* 35(Web Server):W52–W57.
- Guan J, Wang W, Sun B, Fey PD. 2017. Chromosomal targeting by the type III-A CRISPR-Cas system can reshape genomes in *Staphylococcus aureus*. *mSphere.* 2:e00403–17.
- Guo H, Sun S, Eardly B, Finan T, Xu J. 2009. Genome variation in the symbiotic nitrogen-fixing bacterium *Sinorhizobium meliloti*. *Genome* 52(10):862–875.
- Hall JPJ, Brockhurst MA, Harrison E. 2017. Sampling the mobile gene pool: innovation via horizontal gene transfer in bacteria. *Philos Trans R Soc B* 372(1735):20160424.
- Hall JPJ, et al. 2015. Environmentally co-occurring mercury resistance plasmids are genetically and phenotypically diverse and confer variable context-dependent fitness effects. *Environ Microbiol.* 17:5008–5022.
- Hampel S, Cimprich KA. 2016. Conflict resolution in the genome: how transcription and replication make it work. *Cell* 167(6):1455–1467.
- Harrison E, Guymer D, Spiers AJ, Paterson S, Brockhurst MA. 2015. Parallel compensatory evolution stabilizes plasmids across the parasitism-mutualism continuum. *Curr Biol.* 25(15):2034–2039.
- Haug-Baltzell A, Stephens SA, Davey S, Scheidegger CE, Lyons E. 2017. SynMap2 and SynMap3D: web-based whole-genome synteny browsers. *Bioinformatics* 33(14):2197–2198.
- Hesse C, et al. 2018. Genome-based evolutionary history of *Pseudomonas* spp. *Environ Microbiol.* 7:e1002132.
- Holden MTG, et al. 2004. Genomic plasticity of the causative agent of melioidosis, *Burkholderia pseudomallei*. *Proc Natl Acad Sci U S A.* 101(39):14240–14245.
- Holden MTG, et al. 2009. The genome of *Burkholderia cenocepacia* J2315, an epidemic pathogen of cystic fibrosis patients. *J Bacteriol.* 191(1):261–277.
- Hsu PD, Lander ES, Zhang F. 2014. Development and applications of CRISPR-Cas9 for genome engineering. *Cell* 157(6):1262–1278.
- Ito K, Akiyama Y. 2005. Cellular functions, mechanism of action, and regulation of FtsH protease. *Annu Rev Microbiol.* 59: 211–231.
- Jain R, Rivera MC, Lake JA. 1999. Horizontal gene transfer among genomes: the complexity hypothesis. *Proc Natl Acad Sci U S A.* 96(7):3801–3806.
- Janssen PJ, et al. 2010. The complete genome sequence of *Cupriavidus metallidurans* strain CH34, a master survivor in harsh and anthropogenic environments. *PLoS One* 5(5):e10433.
- Kacar B, Garmendia E, Tuncbag N, Andersson DI, Hughes D. 2017. Functional constraints on replacing an essential gene with its ancient and modern homologs. *MBio* 8:e01276–17.
- Kanehisa M, Goto S. 2000. KEGG: Kyoto encyclopedia of genes and genomes. *Nucleic Acids Res.* 28(1):27–30.
- Lercher MJ, Pál C. 2008. Integration of horizontally transferred genes into regulatory interaction networks takes many million years. *Mol Biol Evol.* 25(3):559–567.
- Lillestøl RK, et al. 2009. CRISPR families of the crenarchaeal genus *Sulfolobus*: bidirectional transcription and dynamic properties. *Mol Microbiol.* 72(1):259–272.
- MacLean RC, San Millan A. 2015. Microbial evolution: towards resolving the plasmid paradox. *Curr Biol.* 25(17):R764–R767.
- Makarova KS, et al. 2015. An updated evolutionary classification of CRISPR-Cas systems. *Nat Rev Micro* 13(11):722–736.
- Marraffini LA, Sontheimer EJ. 2010. CRISPR interference: RNA-directed adaptive immunity in bacteria and archaea. *Nat Rev Genet.* 11(3):181–190.
- McGlynn P, Savery NJ, Dillingham MS. 2012. The conflict between DNA replication and transcription. *Mol Microbiol.* 85(1):12–20.
- McInerney JO, McNally A, O'Connell MJ. 2017. Why prokaryotes have pangenomes. *Nat Microbiol.* 2:17040.

- Meinicke P. 2015. UProC: tools for ultra-fast protein domain classification. *Bioinformatics* 31(9):1382–1388.
- Moriya Y, Itoh M, Okuda S, Yoshizawa AC, Kanehisa M. 2007. KAAS: an automatic genome annotation and pathway reconstruction server. *Nucleic Acids Res.* 35(Web Server):W182–W185.
- Nishida H, Abe R, Nagayama T, Yano K. 2012. Genome signature difference between *Deinococcus radiodurans* and *Thermus thermophilus*. *Int J Evol Biol.* 2012: 205274.
- Nowell RW, Green S, Laue BE, Sharp PM. 2014. The extent of genome flux and its role in the differentiation of bacterial lineages. *Genome Biol Evol.* 6:1514–1529.
- Richter M, Rosselló-Móra R. 2009. Shifting the genomic gold standard for the prokaryotic species definition. *Proc Natl Acad Sci U S A.* 106(45):19126–19131.
- Romanchuk A, et al. 2014. Bigger is not always better: transmission and fitness burden of ~1MB *Pseudomonas syringae* megaplasmid pMPPla107. *Plasmid* 73:16–25.
- San Millan A, MacLean RC. 2017. Fitness costs of plasmids: a limit to plasmid transmission. *Microbiol Spectr.* 5(5).
- Shachrai I, Zaslaver A, Alon U, Dekel E. 2010. Cost of unneeded proteins in *E. coli* is reduced after several generations in exponential growth. *Mol Cell* 38(5):758–767.
- Shapiro BJ. 2017. The population genetics of pangenomes. *Nat Microbiol.* 2(12):1574–1574.
- Stern A, Keren L, Wurtzel O, Amitai G, Sorek R. 2010. Self-targeting by CRISPR: gene regulation or autoimmunity? *Trends Genet.* 26(8):335–340.
- Tatusova T, et al. 2016. NCBI prokaryotic genome annotation pipeline. *Nucleic Acid Res.* 44(14): 6614–6624.
- Teeling H, Meyerdierks A, Bauer M, Amann R, Glöckner FO. 2004. Application of tetranucleotide frequencies for the assignment of genomic fragments. *Environ Microbiol.* 6(9):938–947.
- Tett A, et al. 2007. Sequence-based analysis of pQBR103; a representative of a unique, transfer-proficient mega plasmid resident in the microbial community of sugar beet. *ISME J.* 1(4):331–340.
- Tomoyasu T, et al. 1995. *Escherichia coli* FtsH is a membrane-bound, ATP-dependent protease which degrades the heat-shock transcription factor sigma 32. *EMBO J.* 14(11):2551–2560.
- Vercoe RB, et al. 2013. Cytotoxic chromosomal targeting by CRISPR/Cas systems can reshape bacterial genomes and expel or remodel pathogenicity islands. *PLoS Genet.* 9(4):e1003454.
- Vos M, Eyre-Walker A. 2017. Are pangenomes adaptive or not? *Nat Microbiol.* 2(12):1576.
- Wellner A, Gophna U. 2008. Neutrality of foreign complex subunits in an experimental model of lateral gene transfer. *Mol Biol Evol.* 25(9):1835–1840.
- Wick RR, Judd LM, Gorrie CL, Holt KE. 2017. Unicycler: resolving bacterial genome assemblies from short and long sequencing reads. *PLOS Comput Biol.* 13(6):e1005595.
- Wood DE, Salzberg SL. 2014. Kraken: ultrafast metagenomic sequence classification using exact alignments. *Genome Biol.* 15(3):R46.

Associate editor: Tal Dagan